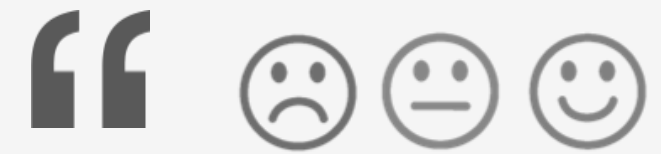


NLP 감성분석 기반
마케팅 시장 분석

NLP 감성 분석에 기반한 온라인 커뮤니티 이슈 키워드 모니터링 시스템

빅데이터미네이터



커뮤니티별 주요 이슈 키워드, 감성 분석 시스템

각 커뮤니티의 제목, 본문, 댓글에 대한 크롤링을 진행하고, 이에 따른 주요 이슈 키워드를 추출한다.

또한 추출한 데이터를 바탕으로 감성 분석하여 결과를 도출하고자 한다.

커뮤니티별 익일 이슈 키워드와 그에 대한 감성 분석 결과를 사용자에게 알려주는 시스템이다.

최종 목표는 매일 온라인 커뮤니티의 주요 이슈 키워드와 그에 대한 감성 분석을
사용자에게 알람으로 전달해주는 서비스이다.

따라서 선정한 커뮤니티인 4곳의
인기 게시판의 하루 치 데이터를 크롤링해 test 데이터를 만들었다.

이를 바탕으로 각 커뮤니티에대한 키워드 추출, 감성 분석 등을 진행해 나갔다.

- 명사 추출, Bag of Words, TF- IDF, 토픽 모델링 등 여러 방식을 조사해보았다. 또한 의미 있는 키워드 추출을 위해 한가지 방식만 적용하는 것이 아닌, 2가지 방식을 결합해서도 시도해 보았다.

결과적으로는 명사 추출을 한 다음 TF-IDF 하는 방식이 프로토타입 단계에서 가장 활용 가능한 결과를 내는 것으로 판단하였다.

```
In [5]: #TF-IDF 이미 한번 명사 단위로 정리 한 명단으로 처리할 시
as_one = ''
for noun in nouns:
    as_one = as_one + ' ' + noun
words = as_one.split()

counts = Counter(words)

vocab = sorted(counts, key=counts.get, reverse=True)

#단어 빈도 정리
word2idx = {word.encode("utf8").decode("utf8"): ii for ii, word in enumerate(vocab,1)}

#역서너리로 정리
idx2word = {ii: word for ii, word in enumerate(vocab)}
```

```
In [6]: #tf-idf
tfidf = TfidfVectorizer(max_features = 10, max_df=0.95, min_df=0)

#generate tf-idf term-document matrix
A_tfidf_sp = tfidf.fit_transform(nouns) #size D x V

In [7]: #tf-idf dictionary
tfidf_dict = tfidf.get_feature_names()
print(tfidf_dict)

morphs_dict = {
    '키워드' : tfidf_dict,
}

df2 = pd.DataFrame(morphs_dict)
df2.to_csv('10키워드s.csv', index=False, encoding="utf-8-sig")

['결혼', '남자', '남편', '본인', '부모', '시간', '엄마', '여자', '자기', '친구']
```

감성 분석: 01. 네이트판

+

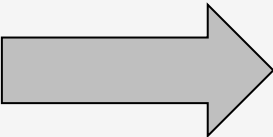
1. 키워드 Top 10 추출

['결혼', '남자', '남편', '본인', '부모', '시간', '엄마', '여자', '자기', '친구']

2. 1위 키워드에 대한 데이터만 추출 및 토큰화

| index | 제목 | 본문 | 댓글 |
|-------|----|----------------------------------------------------------------------|----------------------------------------------------------------------------------------|
| 0 | 0 | 내내음 마라고 하는 시 어머니 어느 순간부터 시 어머니가 친 형 엄마를 지칭해서 할 때 내 내음 마라고 하더라고... | 오이고 진짜 무시 한 여편 내 여편한테 이런 표현 하기 싫지만 대우해 주고 싶지 않... |
| 1 | 1 | 임산에 대한 천정과 시댁의 온도차 추가 여섯달에 세달만 맘에 들 제하고 요즘은 남편 하고 남편 친구하고 토론했던 것이... | 재미도 없고 마지막을 진짜 인생 나옴네 주 착 오지네 임신 초기에 회절대 포복... |
| 2 | 2 | 예비 신랑 할로 때문에 너무 고민인 니다 | 만병하시게 요지는 예랑이와 년 정도 연애 하고 몇 년의 헤어짐을 겪고 결혼 약속을 하... |
| 3 | 3 | 회계사 시절 물더니 헤어지자는 남 자친구 | 근 년간을 남자친구 못바라지하고 많이 증거들었는데 요 저쪽 까라 서로 결혼은 언제 할... |
| 4 | 4 | 원다가 물리친 사람이라 비혼주의가 됐어요 | 부모님을 장례화하고 남자는 알지만 헤어지는 능부시고 어머 나는 물리친 분이세요가난한... |
| 5 | 5 | 시댁 문제로 남편과 심각하게 싸웠어요 | 만병하시게 요 제가 남편이랑 크게 부부싸움을 했는데 게 생각이 잘못한 건지 진짜 의견... |
| 6 | 6 | 예비 사동생 과거 결혼 문제 | 만병하시게 요 년과 연애 후 결혼 진찰 중에 있는 예신이에 고 주 나트 같은 글이예요 라고 말하고 싶은 건가남이랑 본데 집안 트 그렇게 수근 말어지는... |

KNU 감성어 사전
활용하여 감성
분석



| Keyword | content | positive | negative | neutral |
|---------|--------------------------------------------------------|----------|----------|---------|
| 0 | 결혼 ['엄', '사', '어머니', '순간', '사', '어머니', '정', '엄마', ...] | 0 | 0 | 0 |

| | Keyword | negative | neutral | positive |
|---|---------|----------|---------|----------|
| 0 | 결혼 | -397 | 13 | 293 |



⇒ 앞으로 수정해 나가야하는 것

1) '예랑' 등 인터넷 용어와 '온도차', '인성' 등 사전에 없지만 온라인에서 많이 사용되는 감성어를 앞으로 계속해서 사전에 추가해 나가야 한다.

2) 데이터 전처리시, 정확하지 못한 띄어쓰기로 인한 오류를 해결하기 위해, 현재 시도해본, 공백 모두 지우고 pykospacing으로 띄어쓰기 다시 해주는 방식 외에 더 효과적인 방식이 있는지 조사 해볼 것이다.

감성 분석: 02. FM



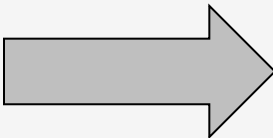
1. 키워드 Top 10 추출

['간부', '군대', '기부', '나라', '백신', '병원', '수술', '아웃', '추진수', '후원']

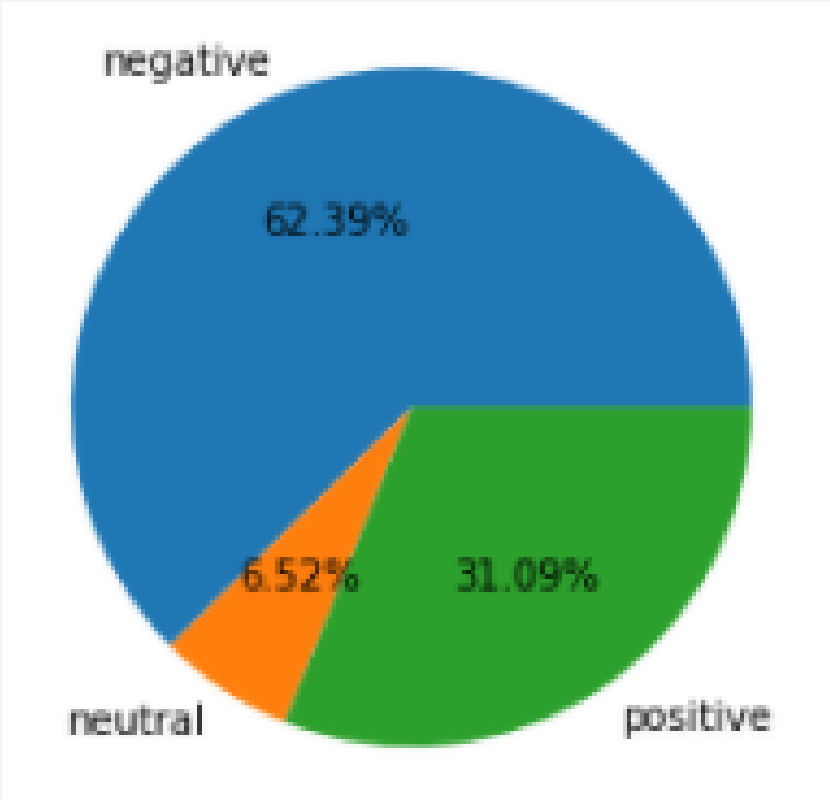
2. 1위 키워드에 대한 데이터만 추출 및 토큰화

'거부', '어머니', '전화', '코로나', '규정', '인하다', '민간', '병원', '수술', '절대', '발다', '전역', '수술', '발다', '하다', '전참', '쓸다', '고통', '견디다', '복무', '계속', '하다', '마비', '통증', '점점', '심해지다', '차례', '수술', '필요하다', '소견', '제출', '하다', '휴가', '요청', '하다', '민간', '병원', '예약', '되다', '휴가', '불가능하다', '까지', '거절', '쓸다', '부모님', '통해', '예약', '하다', '하다', '사진', '만으로는', '예약', '불가능하다', '결국', '휴가', '감수', '병원', '에서', '진료', '발다', '지다', '개월', '지나', '전역', '하다', '추다', '휴가', '나오다', '응급', '수술', '발다', '후유증', '지체', '장애인', '되다', '버리다', '쓸다', '개월', '만원', '재활', '치료', '간병인', '까지', '써다', '되다', '상황', '이지만', '전역', '하다', '이유', '군대', '에서', '하다', '보상', '발다', '있다', '심지어', '간부', '에게', '사과', '조차', '발다', '끝나다', '버리다', '생활', '찾다', '국가', '아들', '다치다', '아들', '군대', '같다', 'ㄴ', '쫓병', '신다', '들다', '살인', '마렵다', '웁다', '군대', '쓰레기', '꼬이다', '임차', '복무', '진급', '문제', '생기다', '되다', '되다', '전역', '하다', '전역', '하다', '책임', '간부', '쏘다', '죽이다', '원래', '민간', '병원', '수술', '치료', '발다', '자다', '보다', '코로나', '핑계', '민간', '수술', '막다', '문제', '간부', '전쟁', '나다', '대가

Polarity 감성
사전 활용하여
감성 분석



| | 부정 | 중립 | 긍정 |
|---|------------|-----------|------------|
| 0 | 244.099291 | 25.490231 | 121.650713 |



⇒ 앞으로 수정해 나가야하는 것

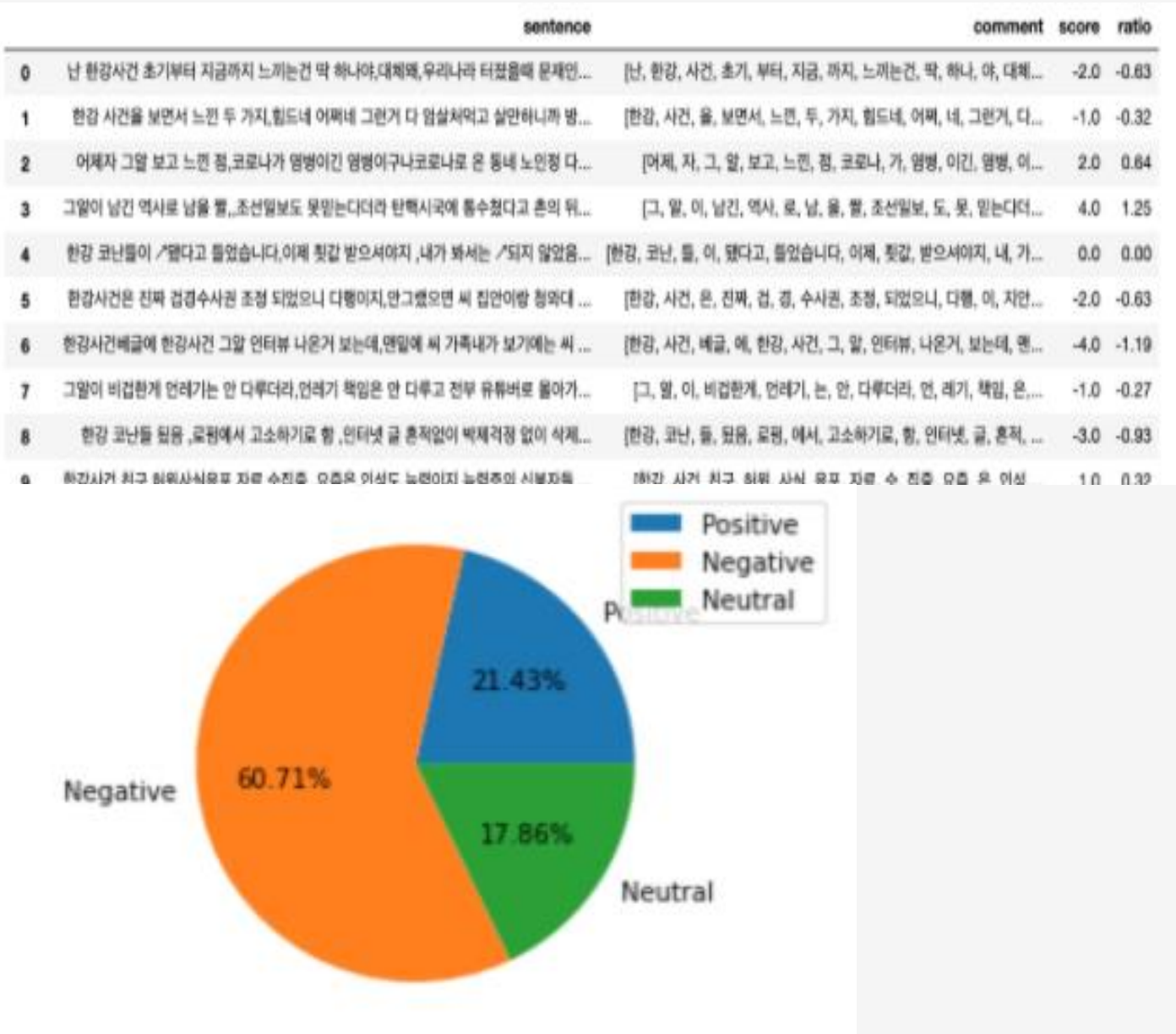
- 1) KNU 감성 사전과 Polarity 감성 사전 중 어느 것이 프로젝트에 더 적합한지 비교할 것이다.
- 2) 토큰화할 때 동사까지 의미 없이 토큰화 되어 버리는 경우가 있어서, 더 나은 토큰화 방법과 데이터 정제 방법을 찾아볼 것이다. 예) 하고싶다 → 하고/싶다

감성 분석: 03. 루리웹

1. 키워드 Top 10 추출

['매수', '민주당', '부모', '사건', '새끼', '유튜브', '의대', '자식', '친구', '한강']

2. “한강“ 키워드 : KNU 감성어 사전 활용하여 감성 분석



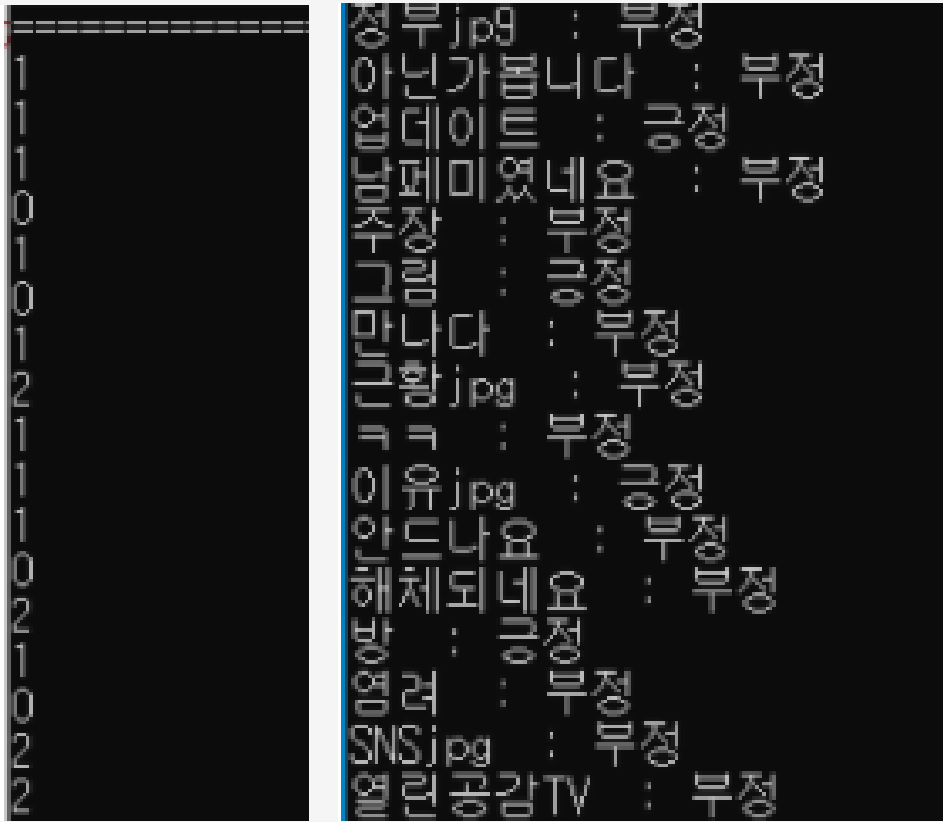
감성 분석: 04. 클리앙

+

1. 키워드 Top 10 추출

```
[('독보적인', 0.5773502691896257), ('포지션의', 0.5773502691896257), ('러면.jpg', 0.5773502691896257)]
[('꼭', 0.0), ('해야', 0.7071067811865475), ('합니다', 0.7071067811865475)]
[('대통령', 0.5505723023275709), ('사저', 0.8347874818836936)]
[('시계', 0.447213595499958), ('보도했던', 0.447213595499958), ('SBS', 0.0), ('이승재', 0.447213595499958), ('기자', 0.447213595499958)]
[('성과가', 0.7071067811865475), ('뭔가요', 0.7071067811865475)]
[('전장관', 0.5773502691896257), ('폐북종하가는', 0.5773502691896257), ('길', 0.0), ('입니다', 0.5773502691896257)]
[('사진관', 0.33563605278348513), ('문재인', 0.3033829386533581), ('대통령', 0.2792927699687521), ('그레고리', 0.4234686473781607), ('추기경을', 0.4234686473781607), ('뵈고', 0.4234686473781607), ('했습니다', 0.4234686473781607)]
[('티트리게', 0.5), ('하는', 0.5), ('중앙', 0.5), ('헛소리', 0.5)]
[('호건', 0.408248290463863), ('매일랜드', 0.408248290463863), ('주지사', 0.408248290463863), ('추모의', 0.0), ('벽', 0.0), ('25만불', 0.408248290463863), ('지적', 0.408248290463863)]
```

2. 감성 분석



최종 발표 전까지

1. 전처리 수정

: 다양한 전처리 방식 계속 시도해보고, 더 성능이 나은 방식 채택

2. 감성어 사전에 단어 추가

: 감성어 사전에 현재 없는 온라인 용어 사전에 계속 추가

3. 전체적인 코드 정리 및 수정

- 감성 분석 방식 채택

- 키워드 10개에 전체에 대한 감성 분석



감사합니다 :)
빅데이터미네이터